













math

Tlan
Report No. 351

510,84

AN ERROR ANALYSIS OF SOLUTIONS TO SPARSE LINEAR PROGRAMMING PROBLEMS

by

Raymond J. Lermit

OCT 2 1 1989

September 17, 1969



DEPARTMENT OF COMPUTER SCIENCE
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN · URBANA, ILLINOIS





AN ERROR ANALYSIS OF SOLUTIONS TO SPARSE LINEAR PROGRAMMING PROBLEMS*

Ву

Raymond J. Lermit

September 17, 1969

Department of Computer Science University of Illinois at Urbana-Champaign Urbana, Illinois 61801

^{*}This work was supported in part by the Advanced Research Projects Agency as administered by the Rome Air Development Center under Contract No. US AF 30(602)4144 and submitted in partial fulfillment of the requirements for the degree of Master of Science in Computer Science, September, 1969.



ABSTRACT

This thesis discusses the round-off errors incurred in solving

Linear Programming problems by the Revised Simplex Method, using the Product

Form of the Inverse - the most usual method of solution. The matrices

involved are assumed to be sparse as is usually the case.

Also discussed are the reinversion of the basis matrix and iterative refinement of the solution obtained.



ACKNOWLEDGEMENT

The author wishes to express his sincere gratitude to his adviser and friend, Professor Robert S. Northcote, for his help in preparing this thesis. He would also like to thank his co-workers Serena Yardley and Teruji Yamamoto.

A special word of thanks to Mrs. Kay Flessner for her patient work in typing a very difficult manuscript.



TABLE OF CONTENTS

			Page
1.	INTRODUCTION		1
2.	NOTATION AND PRELIMINARY RESULTS		3
3.	THE SIMPLEX METHOD	•	10
4.	ERROR ANALYSIS OF THE PRODUCT FORM OF THE INVERSE METHOD		15
5.	REINVERSION OF THE BASIS MATRIX	•	30
6.	ERROR ANALYSIS OF THE REINVERSION ROUTINE		38
7.	ITERATIVE REFINEMENT	•	49
8.	CONCLUDING REMARKS	•	55
LISI	OF REFERENCES	•	57
APPE	ENDIX		58



1. INTRODUCTION

The history of linear programming goes back to the discovery of the simplex algorithm in 1947 by George L. Pantzig, and its implementation on digital computers in the mid 1950's. Little, however, has been written about the round-off errors produced during the computation, the most notable exception being a thesis by Bartels [1].

A linear programming problem may be defined as:

minimize a linear objective function

$$z = c^{T} x$$

subject to the constraints

$$Ax = b, x \geq 0,$$

where A is an m x n real matrix and x, b and c are real vectors of length n, m and n, respectively.

Two approaches to error analysis are possible. The <u>forward error</u> analysis method assumes that a solution $x + \delta x$ has been obtained and the size of δx is then estimated. Alternatively, the <u>backward error analysis</u> approach assumes that the solution obtained is the exact solution of the slightly different problem:

minimize:
$$z + \delta z = (c + \delta c)^{T} x$$

subject to:
$$(A + \delta A) x = b + \delta b, x > 0$$

and estimates of the size of δA , δb and δc are made.

Even the most cursory glance at input data for a linear programming problem will convince the reader that the information is far from exact.

For this reason backward error analysis is a much more realistic approach and is the one which will be pursued in this paper. If a bound can be given on the sizes of δA , δb and δc , and these values are within tolerances determined by inaccuracies in the data, then the problem may be considered solved. This places on the user the responsibility of specifying the accuracy of his data.

A feature of nearly all linear programming problems is that the matrix A is sparse, with very few, often less than one percent, of its elements differing from zero. This fact has been used almost from the outset to reduce the amount of computation and storage required simply by not storing the zeros and omitting calculations involving them. It is also possible to exploit this sparseness in the error analysis.

Control of errors in linear programming has traditionally been done by "tolerances". Quantities that are "very small" are set to zero and numbers greater than some small negative constant are deemed to be positive. Such tolerances are usually set by "the system" but may be changed by the user at his peril.

Efficient solution of problems seems to require these tolerances, and techniques for their automatic control have been developed (e.g., Clasen [2]). However, the final solution should be free of any such artificialities and this is easily achieved by iteratively refining an approximate solution to purify it of its tainted past.

Chapters 3 and 4 cover the simplex method and its backward error analysis respectively. Chapters 5 and 6 cover a reinversion routine that concentrates on retaining sparseness and providing error control. If the solution is not good enough after reinversion the refinement techniques of Chapter 7 may be used. Finally, Chapter 8 shows how slight imperfections in the solution may be removed by perturbing the original problem.

2. NOTATION AND PRELIMINARY RESULTS

It is assumed that all operations are performed in normalized floating-point arithmetic. Define

to be the result of evaluating the expression E. Following Wilkinson [9], the following laws in floating point arithmetic are assumed to hold on a normal machine, providing overflow and underflow do not occur:

fl
$$(a + b) = a(1 + \delta_1) + b(1 + \delta_2)$$

fl $(a - b) = a(1 + \delta_3) - b(1 + \delta_4)$
fl $(a \times b) = (a \times b) (1 + \delta_5)$
fl $(a / b) = (a / b) (1 + \delta_6)$

for any two machine numbers a and b, the quantities $|\delta_i|$, which are relative errors, being bounded by some small constant.

Two degrees of precision are utilized:

normal precision, in which

$$\left|\delta_{i}^{\cdot}\right| \leq \varepsilon_{1}$$

and extended precision where

$$|\delta_i| \leq \epsilon_2$$

No assumptions are made about the relative sizes of ϵ_1 and ϵ_2 . Normal precision is used throughout, except for iterative refinement. Note, however, that many routines use "double precision" for all calculations.

It is assumed that operations involving a zero operand are exact; i.e.,

fl
$$(o + a) = fl (a - o) = a$$

fl $(o - a) = -a$
fl $(o \times a) = fl (a \times o) = 0$
fl $(o / a) = o (a \neq o)$

In practice, however, careful programming should be used to avoid doing these operations at all, as they are redundant.

Upper case letters are used for matrices, lower case letters are used for vectors and for the elements of both matrices and vectors. Where a sequence of matrices or vectors is used, superscripts are employed to distinguish them and their elements are similarly labeled. Matrices and vectors of small terms introduced to account for rounding errors are prefixed by 8 i.e., 8B is a perturbation applied to the matrix B. All vectors are column vectors unless the context implies otherwise.

Throughout the analysis, the ℓ_{∞} vector norm

$$| | v | | = \max_{i} | v_{i} |$$

is used and

$$||A|| = \max ||Ax||$$

$$||x|| = 1$$

$$= \max \sum_{i j=1}^{n} |a_{ij}|$$

for any n \times n matrix A. Since sparse matrices are being used, it is convenient to define the function ∇ by

$$\nabla (t) = \begin{cases} 1 & \text{if } t \neq 0 \\ 0 & \text{if } t = 0 \end{cases}$$

for any scalar t. Thus the number of nonzero elements in a vector

$$v = (v_1, v_2, ..., v_n)$$
 is

$$\sum_{i=1}^{n} \triangledown (v_i)$$

Lemma

If
$$|\theta_k| \le \varepsilon$$
 $k = 1, 2, ..., i$

then

where

$$\epsilon' = \frac{1}{n} \left[\exp \left(\frac{n \epsilon}{1 - \epsilon} \right) - 1 \right]$$
 and $i \le n$ (1)

Proof:

Let
$$p = \prod_{k=1}^{i} (1+\theta_k)$$

$$\geq \prod_{k=1}^{i} (1-\epsilon) = (1-\epsilon)^{i} \qquad (2)$$

$$p \leq (1+\epsilon)^{i}$$

$$p \leq (1 + \epsilon)^{i}$$

$$\leq \left(\frac{1}{1 - \epsilon}\right)^{i} = (1 - \epsilon)^{-i}$$

$$p-1 \leq (1-\epsilon)^{-1}-1 \tag{3}$$

Since

$$-\left[(1-\epsilon)^{i/2}-(1-\epsilon)^{-i/2}\right]^{2} \le 0$$

$$1-(1-\epsilon)^{i} \le (1-\epsilon)^{-i} - 1$$

Hence, from (1),

$$1 - p \le (1 - \epsilon)^{i} - 1 \tag{4}$$

Now (2) and (3) together yield

$$|p-l| \leq (1-\epsilon)^{-1} - l$$

$$= \left(1 + \frac{\epsilon}{1-\epsilon}\right)^{-1} - l$$

$$\leq \frac{i}{n} \left[\exp\left(\frac{n \epsilon}{1-\epsilon}\right) - 1\right]$$

For small \in , \in ' is only slightly larger than \in . It is often convenient to use \in ' instead of \in in error bounds.

In the evaluation of the inner product of two dense vectors, it is necessary to estimate the error in calculating

$$p = fl\begin{pmatrix} n & x_iy_i \\ \sum_{i=1}^n & x_iy_i \end{pmatrix}$$

using floating-point arithmetic. Let

$$\begin{aligned} \mathbf{s_1} &= \mathbf{0} \\ \\ \mathbf{s_{i+1}} &= \mathbf{fl} \ (\mathbf{s_i} + \mathbf{x_i} \mathbf{y_i}) & \mathbf{i} = 1, 2, \dots, n \\ \\ &= \mathbf{s_i} \ (\mathbf{1} + \boldsymbol{\theta_i}) + \mathbf{x_i} \mathbf{y_i} (\mathbf{1} + \boldsymbol{\Delta_i})^2 \\ \\ &\quad \text{where} \quad |\boldsymbol{\theta_i}| \leq \varepsilon, \quad |\boldsymbol{\Delta_i}| \leq \varepsilon \\ \\ &\quad \text{and} \quad \varepsilon = \text{either } \varepsilon_1 \text{ or } \varepsilon_2 \end{aligned}$$

Then

$$p = s_{n+1}$$

$$= \sum_{i=1}^{n} x_i y_i (1 + \Delta_i)^2 \sum_{k=i+1}^{n} (1 + \theta_k)$$

Using the lemma,

Now consider the case where x and y are sparse vectors. Let

$$q_{x} = \sum_{i=1}^{n} \nabla (x_{i})$$

$$q_{y} = \sum_{i=1}^{n} \nabla (y_{i})$$

$$q = \min (q_{x}, q_{y})$$

and

 $q' = \sum_{i=1}^{n} \nabla (x_i y_i)$

Then

$$q' \leq q$$

If x_j , y_j are x_i , y_i , respectively, where $x_i y_i$ is the j-th nonzero term in the sum, then

$$\text{fl} \begin{pmatrix} \sum_{i=1}^{n} & x_{i}y_{i} \end{pmatrix} = \text{fl} \begin{pmatrix} \sum_{j=1}^{q'} & x_{j}y_{j} \end{pmatrix}$$

$$= \sum_{j=1}^{q'} & x_{j}y_{j} & (1 + \Delta_{j}) \\ & & \text{where } |\Delta_{j}| \leq (q' + 1)\epsilon'$$

$$= \sum_{i=1}^{n} & x_{i}y_{i} & (1 + \delta_{i})$$

$$& \text{where } |\delta_{i}| \leq (q + 1)\epsilon'$$

In fact, $\delta_i = \Delta_j$; however, use of the less stringent bound is justified, since q' is hard to estimate but q is easily obtained.

It is sometimes convenient to consider one of the vectors (say y) to be dense, in this case:

$$q_{y} = n$$

$$q = q_x$$

$$fl\begin{pmatrix} n & x_{i}y_{i} \\ i = l & x_{i}y_{i} \end{pmatrix} = \sum_{i=l}^{n} x_{i}y_{i} (l + \delta_{i})$$

where

$$|\delta_i| \le (q_X + 1) \in '$$
.

3. THE SIMPLEX METHOD

A brief outline of a version of the Revised Simplex algorithm is given; a much more detailed analysis may be found in Dantzig [3] or Hadley [6].

The problem to be solved may be put into the form:

minimize

$$z = c^{T}x$$

subject to

$$Ax = b$$

$$x \geq 0$$
.

Let B be the basis matrix $[P_{v_1}, \dots, P_{v_n}]$ where P_{v_1}, \dots, P_{v_n} is a set of m linearly independent columns of A. Then a basic solution is given by:

$$\hat{x} = B^{-1} b$$
,

$$\hat{\mathbf{x}} = [\mathbf{x}_{\mathbf{v}_1}, \dots, \mathbf{x}_{\mathbf{v}_m}]^T,$$

since the inverse exists; this solution is also feasible if $x \ge 0$. The success of the Simplex method depends on the fact that, if an optimal solution exists, there is a basic feasible solution which yields the optimum value of the objective function. It is therefore necessary to consider only the basic feasible solutions in the search for optimality.

An algorithm for solving the problem by the Revised Simplex Method is as follows:

(a) Find an initial basic feasible solution, or indicate that none exists.

(b) Change the basis until optimality is reached, which is achieved by a series of simplex iterations. The steps required for each iteration, starting with some basis $B = [P_{v}, P_{v}, \dots, P_{v}]$ are:

Step 1: Set
$$\gamma^{\mathrm{T}} = [c_{\nu_1}, c_{\nu_2}, \dots, c_{\nu_m}].$$

Step 2: Compute
$$\pi^T = \gamma^T B^{-1}$$
.

Step 3: Compute
$$d_j = c_j - \pi^T P_j$$

for $j = 1$ (1) n, $j \notin \{v_1, v_2, \dots, v_m\}$.

Choose an index s such that

$$d_{s} = \min \{d_{j} : d_{j} < 0\}$$

If no $d_{j} < 0$, the solution is optimal.

Step 5: Find an index r such that

$$\frac{x_{\underline{r}}}{y_{\underline{r}}} = \min_{y_{\underline{i}} > 0} \left\{ \frac{x_{\underline{i}}}{y_{\underline{i}}} \right\}$$

If no y_i is > 0, then the solution is unbounded Choosing the minimum of x_i/y_i assures that the new basic solution x' will be feasible.

Step 6: Change the basis to

$$B' = [P_{v_1}, ..., P_{v_{r-1}}, P_s, P_{v_{r+1}}, ..., P_{v_m}].$$

Step 7: Compute the new solution,
$$x'$$
, using $x' = (B')^{-1} b$.

Since three sets of equations are solved using the matrix B, it is advantageous to have B^{-1} available. If B^{-1} is known, then the inverse of the new basis $(B')^{-1}$ is easily obtained, since

$$(B')^{-1} = (B')^{-1} B B^{-1}$$

= $(B^{-1} B')^{-1} B^{-1}$.

B' is identical to B except that the r-th column, P_{v_r} , is replaced by P_s , so B⁻¹ B' takes on the simple form:

where
$$y = [y_1, y_2, \dots, y_m]^T = B^{-1} P_s$$

Such a matrix is called an <u>elementary column matrix</u>, or simply an <u>elementary matrix</u>. It is easily inverted to give:

$$\left(\mathbb{B}^{-1} \mathbb{B}^{1} \right)^{-1} = \begin{bmatrix} 1 & -y_{1}/y_{r} \\ 1 & -y_{2}/y_{r} \\ & \vdots \\ & 1/y_{r} \\ & \vdots \\ & -y_{m}/y_{r} \\ & 1 \end{bmatrix} = T$$

say, and

$$(B')^{-1} = T B^{-1}$$

Clearly, if the starting basis is B_0 , and T^1 , T^2 , ..., T^p are the elementary matrices produced during p iterations, then the inverse of the current basis B_p , say, is given by:

$$B_{p}^{-1} = T^{p} T^{p-1} \dots T^{2} T^{1} B_{0}^{-1}$$

It is usual to take $B_0 = B_0^{-1} = I$, so that

$$B_p^{-1} = T^p T^{p-1} \dots T^2 T^1.$$

If a unit matrix cannot be found from the columns of A, then one must be manufactured using "artificial" columns. When the corresponding artificial variables are all removed (or reduced to zero) by subsequent iterations, a basic feasible solution has been found.

When the matrix B_p^{-1} is not calculated explicitly, but only T^p , T^{p-1} , ..., T^2 , T^1 are stored, the method is called the <u>product form of the inverse</u>. We will restrict our analysis to this method, since it is used more extensively than any other method in large scale LP algorithms.

In order to reduce round off errors, storage requirements and running time, it is desirable from time to time to start afresh and recalculate the inverse of the current basis. This technique, called reinversion, is discussed in Chapter 5.

4. ERROR ANALYSIS OF THE PRODUCT FORM OF THE INVERSE METHOD

In this chapter the errors introduced during one iteration of the Simplex method are considered.

Suppose that, after p iterations, the approximate inverse of the basis B is given in product form by:

where each T^{j} is an elementary matrix. Define an error matrix E by:

$$T^{p} T^{p-1} \dots T^{2} T^{1} (B + E) = I.$$
 (1)

Suppose also that, in performing one iteration, the r-th column of B is replaced by a non-basis column of A, $A_s = w$ say, giving a new basis matrix B'. A new elementary matrix T^{p+1} is formed which satisfies:

$$T^{p+1} T^p T^{p-1} \dots T^2 T^1 (B^i + E^i) = I$$
 (2)

The round off errors in (1) and (2) are measured by ||E|| and $||E^t||$, respectively. The problem is to find a bound for $||E^t||$ in terms of ||E||.

To simplify the notation assume, without loss of generality, that the j-th elementary matrix, T^j , is a unit matrix except for its j-th column and that in the j-th iteration, it is the j-th column of B that is replaced. Thus, pivoting occurs down the diagonal. This can be achieved by suitable interchange of rows and columns. In practice the actual pivots may be recorded by saving the pivot column number corresponding to each component of the basic solution. Thus set the j-th elementary matrix to a unit matrix except for the j-th column which is

$$[\eta_1^j, \eta_2^j, \ldots, \eta_m^j]^T$$

Since T^j is obtained by inverting an elementary matrix whose j-th column is y^j say, define the vector η^j to be given exactly by

$$\eta_{i}^{j} = \begin{cases} -y_{i}^{j} / y_{j}^{j} & i = 1(1)m, i \neq j \\ 1 / y_{j}^{j} & i = j. \end{cases}$$

Whenever the T^{j} is used in calculation, the error introduced by the divisions will be included in the error analysis.

Either of two different methods may be used in calculating η^{j} , in each of which $t=-y^{j}_{j}$ and $y^{j}_{j}=-1$ are set first. One method is then to divide each element of y^{j} by t. The other is to multiply each element by t^{-1} . The latter method will be faster if division takes longer than multiplication; it does, however, produce two round off errors for each element. It is assumed that the latter method is used.

One step in the Simplex algorithm is to compute

$$v = B^{-1} w$$

for which the approximate solution is given by

$$y = fl (T^p T^{p-1} \dots T^2 T^1 w)$$
 (3)

Setting

$$\alpha^1 = w$$

and calculating successively

$$\alpha^{k+1} = fl(T^k \alpha^k)$$
 (k = 1, 2, ..., p) (4)

yields

$$y = \alpha^{p+1}$$
.

Equation (4) is equivalent to

$$\alpha_{i}^{k+1} = \begin{cases} \text{fl} (\alpha_{i}^{k} + \eta_{i}^{k} \alpha_{k}^{k}) & \text{i=l, 2, ..., m} \\ \text{fl} (\eta_{k}^{k} \alpha_{k}^{k}) & \text{i=k} \end{cases}$$
 (5)

Rearrangement of (5) to include the effect of computations on previous $\alpha^{\mathbf{k}}$ gives:

$$\alpha_{i}^{\ell} = \begin{cases} \text{fl} \left(\sum_{k=i}^{\ell-1} \alpha_{k}^{k} \eta_{i}^{k}\right) &, & i < \ell \end{cases} (6)$$

$$\text{fl} \left(w_{i} + \sum_{k=1}^{\ell-1} \alpha_{k}^{k} \eta_{i}^{k}\right) &, & i \geq \ell \end{cases}$$

In particular, setting $i = \ell$, yields

For convenience in notation, let $\beta_k = \alpha_k^k$. Then

and

$$y_{i} = \alpha_{i}^{p+1} = \begin{cases} \text{fl } (\sum_{k=i}^{p} \beta_{k} \eta_{i}^{k}) & i \leq p \\ p & p \\ \text{fl } (w_{i} + \sum_{k=1}^{p} \beta_{k} \eta_{i}^{k}) & i > p \end{cases}$$
 (9)

Now let L be the matrix

$$L = \begin{bmatrix} 1 & & & & & \\ -\eta_{2}^{1} & 1 & & & & \\ -\eta_{3}^{1} & -\eta_{3}^{2} & 1 & & & \\ \vdots & \vdots & & \ddots & & \\ -\eta_{m}^{1} & -\eta_{m}^{2} & \dots & 1 \end{bmatrix}$$

Then (8) may be written in the form:

$$L \beta = w$$

from which β is uniquely given by

$$\beta = L^{-1} w$$
.

Using the relations given in Lemma 1 of the Appendix, it can be seen that β exactly satisfies

$$(L + \delta L) \beta = w + \delta w$$
,

$$(s L)_{ij} = \begin{cases} \eta_i^j & s_{ij} \\ 0 & i \leq j \end{cases},$$

$$|\delta_{ij}| \leq (q_i + 3) \epsilon_1^{i}$$

$$q_i = \sum_{k=1}^{i-1} \nabla (\eta_i^k)$$

$$|\delta w_i| \leq q_i |w_i| \in$$

If
$$\Delta w = L \beta - w$$

then $\triangle w = \delta w - \delta L \cdot \beta$

$$||\Delta w|| \le ||\delta w|| + ||\delta L|| ||\beta||$$

$$\le q_{L} ||w|| \epsilon_{1}^{s} + ||\delta L|| ||\beta||$$
(10)

where $q_L = \max_{i} q_{i}$

then
$$|\delta L| \leq \epsilon_1^* q_L (q_L + 3) M_L$$
 (12)

and hence

$$||\Delta w|| \le \epsilon_{1}^{*} (q_{L}^{||w||} + q_{L}^{*} (q_{L}^{*} + 3) M_{L}^{||\beta||})$$
 (13)

Now estimate the error involved in calculating y using (9).

Let U be the matrix

provided that $p \leq m$. Then

$$y = fl(U\gamma)$$

where
$$\gamma = [\beta_1, \beta_2, \dots, \beta_p, w_{p+1}, \dots, w_m]^T$$
 (14)

By using the results of Lemma 2 in the Appendix, it can be shown that y satisfies

$$y = U (\gamma + \delta \gamma)$$

where

$$\delta \gamma = [\delta \beta_1, \delta \beta_2, \ldots, \delta \beta_p, \Delta w_{p+1}, \ldots, \Delta w_m]^T$$
 say,

$$||s\gamma|| \leq \frac{||u^{-1}||}{1 - 2 \epsilon_{1}' q_{U}(q_{U}+2) M_{U} ||u^{-1}||} \times \left[\epsilon_{1}' q_{U}(q_{U}+2) M_{U} ||y|| + \epsilon_{1}' ||\gamma||\right]$$

$$q_{U} = \max_{i} \begin{cases} \sum_{k=1}^{p} \nabla (\eta_{i}^{k}) & , i \geq p \\ \sum_{k=1}^{p} \nabla (\eta_{i}^{k}) + 1 & , i
$$(15)$$$$

= the maximum number of nonzero elements in any row of U.

and

$$\begin{aligned} & \mathbf{M}_{\mathbf{U}} &= \max_{\mathbf{i},\mathbf{j}} \; (|\mathbf{U}_{\mathbf{i}\mathbf{j}}|). \end{aligned}$$
 If $\delta \beta = (\delta \beta_{\mathbf{l}}, \delta \beta_{\mathbf{l}}, \ldots, \delta \beta_{\mathbf{p}})$ and $\Delta \mathbf{w}^{\mathbf{i}} = (\Delta \mathbf{w}_{\mathbf{p}+\mathbf{l}}, \ldots, \Delta \mathbf{w}_{\mathbf{n}})$
$$||\delta \beta|| \leq ||\delta \gamma|| \text{ and } ||\Delta \mathbf{w}^{\mathbf{i}}|| \leq ||\delta \gamma||$$

then

The perturbation $\delta \beta$ of β gives rise to a perturbation L $\delta \beta$ in w_1 , ..., w_p Therefore, combining (13) and (15)

$$y = fl (T^p T^{p-1} ... T^2 T^1 w)$$

= $T^p T^{p-1} ... T^2 T^1 (w + v)$

where
$$||v|| \le \max \{ \in (q_L ||w|| + q_L (q_L + 3) M_L ||\beta||) + |L|||8\beta||, ||8\gamma|| \}$$
 (16)

and
$$||L|| \leq q_L M_L$$

To simplify (16), let

$$M = \max_{i,k} (|\eta_i^k|, 1)$$
 (17)

Then
$$M_{L} \leq M, M_{U} \leq M.$$
 (18)

Let
$$q = \max_{i} \sum_{k=1}^{p} \nabla (\eta_{i}^{k})$$

Then
$$||v|| \le \epsilon_{1}^{i} (q||w|| + q (q+3) M ||\beta||)$$

 $+ \frac{\epsilon_{1}^{i} q M || U^{-1}||}{1 - 2 \epsilon_{1}^{i} q (q+3) M ||U^{-1}||} [q (q+3) M ||y|| + ||\beta|| + q||w||]$ (19)

The only quantity in the above inequality which is unknown is $||U^{-1}||$; fortunately approximations to the elements of U^{-1} are found during the calculation of the Π vectors.

Consider now the calculation of the η vectors forming the approximate inverse of β , and consider β to be a unit matrix except for the first p columns. If

$$B = \begin{bmatrix} \alpha_{11}^1 & \alpha_{12}^1 & \dots & \alpha_{1p}^1 \\ \alpha_{21}^1 & \alpha_{22}^1 & \dots & \alpha_{2p}^1 \\ \vdots & & & \vdots \\ \alpha_{p1}^1 & \alpha_{p2}^1 & \dots & \alpha_{pp}^1 \\ \vdots & & & \vdots \\ \alpha_{m1}^1 & \alpha_{m2}^1 & \dots & \alpha_{mp}^1 & \dots \end{bmatrix}$$

and

$$\alpha_{i,j}^{k+1} = \begin{cases} \text{fl } (\alpha_{i,j}^{k} + \alpha_{k,j}^{k} \, \eta_{i}^{k}) & \text{i } \neq k \\ \text{fl } (\alpha_{k,j}^{k} \, \eta_{k}^{k}) & \text{i } = k \end{cases}$$
for $i = 1(1)$ m , $j = 1(1)$ p , $k = 1(1)$ j-1

then

$$\alpha_{ij}^{j} = \begin{cases} \text{fl} \left(\sum_{k=i}^{j-1} \alpha_{kj}^{k} \eta_{i}^{k} \right) & i < j \\ \text{fl} \left(\alpha_{ij}^{1} + \sum_{k=1}^{j-1} \alpha_{kj}^{k} \eta_{i}^{k} \right) & i \geq j \end{cases}$$

$$(21)$$

and the Π vectors are given by

$$\eta_{i}^{j} = \begin{cases}
\text{fl} \left(-\alpha_{i,j}^{j}/\alpha_{j,j}^{j}\right) & \text{i } \neq \text{j} \\
\text{fl} \left(1/\alpha_{j,j}^{j}\right) & \text{i } = \text{j}
\end{cases}$$
(22)

Let \(\text{be the matrix} \)

Then

$$(UT)_{i,j} = \begin{cases} \sum_{k=i}^{j} \alpha_{k,j}^{k} \eta_{i}^{k} & i \leq p, j \leq p \\ \sum_{k=1}^{i} \alpha_{k,j}^{k} \eta_{i}^{k} + \alpha_{i,j}^{l} & i > p, j \leq p \\ 1 & i = j > p \\ 0 & \text{otherwise} \end{cases}$$

But, for $i < j \le p$,

$$\sum_{k=i}^{j} \alpha_{kj}^{k} \eta_{i}^{k} = \alpha_{jj}^{j} \eta_{i}^{j} + \sum_{k=i}^{j-1} \alpha_{kj}^{k} \eta_{i}^{k}$$

and, from (21),

$$\alpha_{\mathbf{i},\mathbf{j}}^{\mathbf{j}} = \sum_{k=\mathbf{i}}^{\mathbf{j}-\mathbf{l}} \alpha_{\mathbf{k},\mathbf{j}}^{\mathbf{k}} \ \eta_{\mathbf{i}}^{\mathbf{k}} + \sum_{k=\mathbf{i}}^{\mathbf{j}-\mathbf{l}} \mathbf{s}_{\mathbf{i},\mathbf{j}}^{\mathbf{k}} \ \alpha_{\mathbf{k},\mathbf{j}}^{\mathbf{k}} \ \eta_{\mathbf{i}}^{\mathbf{k}}$$

$$\mid \delta_{i,j}^{k} \mid \leq (q_i + 2) \in 1$$

$$q_i = \sum_{k=1}^{p} \nabla (\eta_i^k)$$

$$\sum_{k=i}^{j} \alpha_{kj}^{k} \eta_{i}^{k} = \alpha_{ij}^{j} + \alpha_{jj}^{j} \eta_{i}^{j} + \sum_{k=i}^{j-1} \delta_{ij}^{k} \alpha_{kj}^{k} \eta_{i}^{k} .$$
 (24)

Similarly, for i > p,

$$\sum_{k=1}^{i} \alpha_{kj}^{k} \eta_{i}^{k} + \alpha_{ij}^{l} = \alpha_{ij}^{j} + \alpha_{jj}^{j} \eta_{i}^{j} + \sum_{k=1}^{j-1} \delta_{ij}^{k} \alpha_{kj}^{k} \eta_{i}^{k} + \Delta_{ij} \alpha_{ij}^{l} \alpha_{ij}^{l}$$

$$+ \Delta_{ij} \alpha_{ij}^{l}$$

$$(25)$$

but, from (22),

$$\alpha_{ij}^{j} + \alpha_{jj}^{j} \eta_{i}^{j} = 2 \delta_{ij}^{i} \alpha_{ij}^{j} \qquad i < j$$

$$\alpha_{jj}^{j} \eta_{j}^{j} \qquad = 1 + \alpha_{jj}^{j} \eta_{j}^{j} \delta_{jj}^{j}$$

$$(26)$$

where

and

$$|\delta_{i,j}| \leq \epsilon_{i}$$

Hence

$$U\Gamma = I + F$$

where F is a matrix of small terms caused by round off errors.

If
$$q = \max_{i} \sum_{k=1}^{p} \nabla (\eta_{i}^{k})$$

$$M = \max_{i} (|\eta_{i}^{k}|, 1),$$

as before, G is the matrix

and H is a $p \times p$ matrix given by

$$\text{Then} \qquad \begin{cases} \alpha_{\mathbf{i}\mathbf{j}}^{\mathbf{j}} & \eta_{\mathbf{j}}^{\mathbf{j}} & \delta_{\mathbf{j}\mathbf{j}}^{\mathbf{j}} \\ \alpha_{\mathbf{k}\mathbf{j}}^{\mathbf{j}} & \eta_{\mathbf{j}}^{\mathbf{k}} & \delta_{\mathbf{i}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} \\ \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf{k}} & \lambda_{\mathbf{k}\mathbf{j}}^{\mathbf$$

Hence

$$\begin{split} \sum_{j=1}^{m} |F_{i,j}| &\leq \sum_{j=1}^{p} (|\sum_{k=1}^{j} \delta_{i,j}^{k} \alpha_{k,j}^{k}|^{2} |\eta_{i}^{k}|) \\ &+ 2\varepsilon_{i}^{*} \sum_{j=1}^{p} |\alpha_{i,j}^{j}| + \varepsilon_{i}^{*} \sum_{j=1}^{p} |\alpha_{i,j}^{l}| \\ &\leq \varepsilon_{i}^{*} [q(q+2)M||G|| + 2||H|| + ||B||] \end{split}$$

Thus

$$||F|| \le \epsilon_1' [q (q+2) M ||G|| + 2 ||H|| + ||B||]$$
 (27)

and

$$||\mathbf{r}|| \le \max \{ ||\mathbf{G}||, ||\mathbf{B}|| + 1 \}$$

Now

$$U^{-1} = \Gamma (I + F)^{-1}$$

Substituting the bound for $||U^{-1}||$ into equation (19) gives a bound for ||v||.

It has been shown that

$$y = fl (T^p T^{p-1} ... T^2 T^l, w)$$

$$= T^p T^{p-1} ... T^2 T^l (w + v).$$

It is now possible to estimate | | E' | |, where E' satisfies:

$$T^{p+1} T^{p} . . . T^{2} T^{1} (B' + E') = I.$$

The new elementary matrix \textbf{T}^{p+1} has as its p+1 -st column the vector $\boldsymbol{\eta}^{p+1}$ where:

$$\eta_{i}^{p+1} = \begin{cases}
-y_{i} / y_{p+1} & i \neq p+1 \\
1 / y_{p+1} & i = p+1
\end{cases}$$

Thus

$$T^{p+1}$$
 $y = [0, \dots, 0, 1, 0, \dots, 0]^T$

$$= T^{p+1} T^p \dots T^2 T^1 (w + v)$$

since w is the p + 1 -st column of B', v is the p+1 - st column of E'.

For all other columns B' and B are identical, and, since

$$T^{p+1} \cdot \cdot \cdot T^2 T^1 (B + E) = T^{p+1} I$$

$$= T^{p+1},$$

(B' + E') and (B + E) are the same except for column p+1; i.e., E' is E with the p+1 -st column replaced by v. Hence

$$||E'|| \le ||E|| + ||v||$$

Finding a bound for ||v|| is not as difficult as equations (19) and (27) would seem to indicate since the only quantities required that are not recorded are the elements of G. (β is just the last column of G). These values are obtained during the calculation of y and may be used to update ||G|| from iteration to iteration.

This bound depends heavily on:

- 1. The sparseness of the η vectors;
- 2. The size of their elements.

It is, therefore, very important to reduce the magnitude of the η elements by controlling the size of 'pivots'. Consideration also should be given to maintaining spareseness.

5. REINVERSION OF THE BASIS MATRIX

During the execution of the simplex algorithm, the choice of pivot is restricted at each iteration by the need to obtain a significant improvement in the value of the objective function and to maintain feasibility. Little can be done to control round-off errors and the sparseness of the inverse. In recalculating the inverse a flexible pivot selection agenda can be adopted, and the errors and density be thereby greatly reduced.

Before discussing the general reinversion algorithm, the decomposition of a sparse matrix A by two different methods will be considered. In both cases pivoting occurs down the diagonal.

(1) Product form decomposition.

$$A = E^{1}E^{2} \dots E^{n-1}E^{n}$$

where

(2) Decomposing A into a product of lower and upper triangular matrices

$$L = \begin{bmatrix} \ell_{11} & & & & \\ \ell_{21} & \ell_{22} & & & \\ \ell_{31} & \ell_{32} & \ell_{33} & & \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \ell_{n1} & \ell_{n2} & \ell_{n3} & \cdots & \ell_{nn} \end{bmatrix}$$

and

so that A = LU.

Having decomposed the matrix, inversion is easy, since

$$(E^{j})^{-1} = \begin{bmatrix} 1 & & -e_{1j}/e_{jj} & & & & \\ & 1 & & \vdots & & & \\ & & \cdot & \vdots & & & \\ & & & 1/e_{jj} & 1 & & \\ & & & -e_{nj}/e_{jj} & & \cdot & 1 \end{bmatrix}.$$

$$L^{-1} = L^{n} L^{n-1} \cdot \cdot \cdot L^{2} L^{1}$$

$$\mathbf{L}^{\mathbf{j}} = \begin{bmatrix} \mathbf{1} & \mathbf{i} & & & & \\ & \mathbf{1} & & & & \\ & & & \mathbf{1} & & \\ & & & \mathbf{1}/\ell_{\mathbf{j}\mathbf{j}} & & \\ & & & -\ell_{\mathbf{j}+\mathbf{1}\mathbf{j}}/\ell_{\mathbf{j}\mathbf{j}} & & \mathbf{1} & \\ & & & & -\ell_{\mathbf{n}\mathbf{j}}/\ell_{\mathbf{j}\mathbf{j}} & & & \mathbf{1} \end{bmatrix}$$

and similarly for U^{-1} . Note that the pattern of zeros and nonzeros in E^{j} , L^{j} , U^{j} correspond exactly to that of the matrices from which they were derived.

The density of the product form decomposition is at least as great as in the LU form, as can be shown thus:

$$= S^{j} T^{j} (say).$$

Then
$$A = S^1 T^1 S^2 T^2 \dots S^{n-1} T^{n-1} S^n T^n$$
.
For $i \le j$, $T^i S^j = S^j T^i$

So the product of the S's and T's may be rearranged to give:

$$A = S^{1} S^{2} \dots S^{n-1} S^{n} T^{1} T^{2} \dots T^{n-1} T^{n}$$

 s^1 , s^2 , ..., s^n are all lower triangular matrices, so $L = s^1 s^2 \dots s^n$ is lower triangular.

Similarly $U=T^1\ T^2\dots T^n$ is upper triangular, and LU is the (unique) decomposition into lower and upper triangular parts.

$$L = S^1 S^2 \dots S^n$$

$$= \begin{bmatrix} e_{11} & & & & & \\ e_{21} & 1 & & & & \\ \vdots & \ddots & & e_{23} & 1 & \\ e_{n1} & & 1 & e_{n2} & & 1 \end{bmatrix} \cdots \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & & \ddots & \\ & & & e_{n-1,n-1} & & \\ & & & & e_{n-1,n-1} & & \\ & & & & & e_{n-1,1} & & \\ & & & & & e_{nn} \end{bmatrix}$$

Hence L contains the same number of nonzero elements as in S^1 , S^2 , ..., S^n (excluding unit diagonals which need not be stored).

$$U = T^{1} T^{2} \dots T^{n}$$

$$U = (T^{n})^{-1} (T^{n-1})^{-1} \dots (T^{2})^{-1} (T^{1})^{-1}$$

$$= \begin{bmatrix} 1 & & & -e_{1n} \\ & 1 & & & -e_{2n} \\ & & & \ddots & \vdots \\ & & & & -e_{n-1n} \\ & & & & 1 \end{bmatrix} \begin{bmatrix} 1 & & & -e_{1n-1} \\ & 1 & & & -e_{2n-1} \\ & & & & \ddots & \vdots \\ & & & & & 1 \end{bmatrix} \cdots \begin{bmatrix} 1 & -e_{12} \\ & 1 \\ & & & & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -e_{12} & -e_{13} & \dots & -e_{1n} \\ & 1 & -e_{23} & \dots & -e_{2n} \\ & & 1 & -e_{23} & \dots & -e_{2n} \\ & & & 1 & \dots & -e_{3n} \\ & & & & \ddots & \vdots \\ & & & & & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -e_{12} & -e_{13} & \dots & -e_{1n} \\ & 1 & -e_{23} & \dots & -e_{2n} \\ & & & & \ddots & \vdots \\ & & & & & \ddots & \vdots \\ & & & & & & \ddots & \vdots \\ & & & & & & \ddots & \vdots \\ & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & \ddots & \vdots \\ & & & & & & & & \ddots & \vdots \\ & & & & & & & & \ddots & \vdots \\ & & & & & & & & & \ddots & \vdots \\ & & & & & & & & & & & & & & \\ & & & & & & & & & & & & \\ & & & & & & & & & & & \\ & & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & &$$

Therefore the number of nonzero elements in T^1 , T^2 , . . . , T^n is equal to that in U^{-1} (again excluding the unit diagonals).

Since the density of U⁻¹ will be greater than that of U, due to "fill in" during the inversion procedure, it may be concluded that the decomposition into elementary matrices will, in general, give rise to more non-zero elements than will LU decomposition. (This ignores the possibility of cancellation, which might produce zero elements from sums of nonzero ones; however, this should not often happen in practice).

It is therefore desirable to use an LU decomposition method for reinversion, not only to reduce the rounding errors, but also to reduce the number of Π elements representing the inverse.

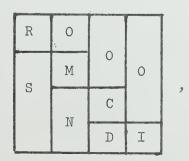
It has been shown that a triangular matrix decomposes readily into elementary matrices. It is therefore desirable to sort the rows and columns of the tasis to form a matrix which is as close to a triangular matrix as possible. The sorted basis is then partitioned as follows:

	R	M	C	N
R	$\begin{bmatrix} x & & & \\ x & x & \\ x & x & x \\ x & x &$	0	0	0
М	x x x x x x x x x x x x x x x x x x x	X X X X X X X X X X X X X X X X X X X	0	0
С	X	X	$\begin{array}{c} x \\ x $	0
N	X	x x x x x x	X	

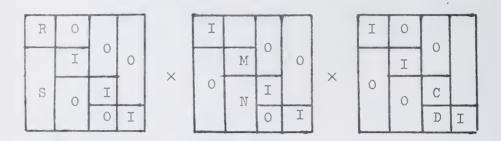
In practice the matrix need not be physically sorted, all that is needed are the indices of the pivot positions.

The NN block contains logical vectors in the basis; these will produce no η records. The RR and CC blocks are triangular. The MM block contains rows and columns which cannot be triangularized by sorting. LU decomposition is applied on this block.

Rewrite the matrix, in partitioned form, as:

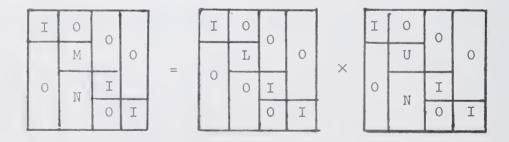


where O and I are zero and identity matrices of appropriate sizes. This decomposes to:



The first and last matrices are lower triangular, pivoting down the diagonal gives Π vectors with no increase in the number of nonzero elements.

 $\label{eq:lower_state} \mbox{If } \mbox{M} = \mbox{LU, L and U being lower and upper triangular respectively,} \\ \mbox{then}$



The first matrix is lower triangular and the second can be transformed by row and column interchanges into:

I	0		0
	I	0	N
0	0	I	
		0	U

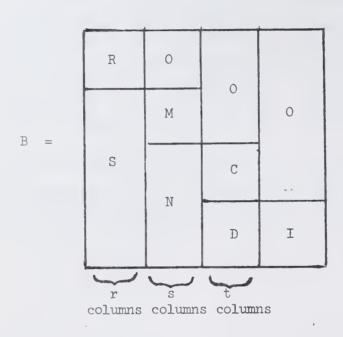
which is upper triangular. After inversion, the interchanges can be reversed. The only place where elements are added to the inverse is in the formation of L and U. Since the pivoting strategy within the M block may be employed, this may be used to minimize the number of nonzeros created.

The exact form of the pivoting strategy depends on such factors as the number of columns of the matrix that can be stored in memory at one time. An important point is to chose the pivots dynamically, at each step choosing that pivot which minimizes the number of nonzeros created.

It may be necessary to change the pivot selection method where choosing a particular pivot would result in very large elements in the representation of the inverse. Discussion of this point is deferred until the next chapter.

6. ERROR ANALYSIS OF THE REINVERSION ROUTINE

In this chapter the computation of the inverse of the matrix



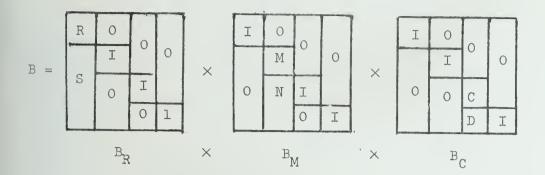
using the product form

$$T^{p} T^{p-1} \dots T^{2} T^{1}, p = r + s + t,$$

is considered. It is shown that

$$T^p T^{p-1} \dots T^2 T^1 (B + \delta B) = I$$

i.e., T^p . . . T^l is the exact inverse of some perturbed matrix B + δB , and a bound for $||\delta B||$ is established. As before B is split into the form



Consider first the inversion of

The η vector for the j^{th} elementary matrix T^j (j=l(l) r) is given by:

$$\eta_{j}^{j} = \text{fl } (1/b_{jj}) = \frac{1}{b_{jj}(1+\delta_{jj})}$$

$$\eta_{i}^{j} = \text{fl } (-b_{ij} \times \eta_{j}^{j}) = \frac{-b_{ij}(1+\delta_{ij})}{b_{jj}(1+\delta_{jj})} \qquad \text{for } i > j$$

where
$$|\delta_{ij}| \leq \epsilon_1$$
,

The terms b_{ij} $(1 + \delta_{ij})$ may be regarded as the elements of a matrix; hence T^r . . . T^l is the exact inverse of the perturbed matrix $(B_R + \delta B_R)$ where

$$\left(\delta B_{R} \right)_{ij} = \begin{cases} \delta_{ij} & \left(B_{R} \right)_{ij} \\ \vdots & \vdots \\ 0 & \text{otherwise} \end{cases}$$

Similarly, the product of the elementary matrices produced for the inversion of B_{C} is the exact inverse of $(B_{C}+\delta B_{C})$ where

$$(\delta B_C)_{ij} = \begin{cases} \delta_{ij} (B_C)_{ij} & j = r + s + l(1)p \\ i = j(1)m \\ 0 & \text{otherwise}$$

and
$$|\delta_{i,j}| \leq \epsilon_1$$

This leaves only the matrix

$$B_{M} = \begin{bmatrix} I & O & & & \\ & M & & O \\ & & M & & O \\ & & I & & \\ & & O & I \end{bmatrix}$$

to be inverted. Since the elements of N are simply copied into the η vectors with a change in sign, the only error is that in inverting M. Therefore set,

$$\mathbf{M} = \begin{bmatrix} \mathbf{m}_{11}^1 & \mathbf{m}_{12}^1 & \cdots & \mathbf{m}_{1s}^1 \\ \mathbf{m}_{21}^1 & \mathbf{m}_{22}^1 & \cdots & \mathbf{m}_{2s}^1 \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & & \\ & &$$

and compute a sequence of matrices M^1 , M^2 , ..., M^s , M^{s+1} , by forming a new N vector from the k-th column of M^k using only elements on or below the diagonal.

Thus

$$\begin{split} & \boldsymbol{\eta}_k^k = \text{fl } (1/\boldsymbol{m}_{kk}^k) = \underbrace{\frac{1}{\boldsymbol{m}_{kk}^k (1 + \boldsymbol{\delta}_{kk}^k)}} \\ & \boldsymbol{\eta}_i^k = \text{fl } (-\boldsymbol{m}_{ik}^k \boldsymbol{\eta}_k^k) = \underbrace{\frac{-\boldsymbol{m}_{ik}^k (1 + \boldsymbol{\delta}_{ik}^k)}{\boldsymbol{m}_{kk}^k (1 + \boldsymbol{\delta}_{kk}^k)}} \end{split}$$

where
$$\left|\delta_{jk}^{k}\right| \leq \epsilon_{j}^{\prime}$$
. (1)

The matrix M^k is then updated by this new elementary matrix giving M^{k+1} . If there were no rounding errors, the elements of M^{k+1} would be zero in the k-th column below the diagonal and M^{k+1}_{kk} would be unity. If we set these to zero and one respectively then:

$$\mathbf{m}_{i,j}^{k+1} = \begin{cases} 1 & i = j = k \\ 0 & i > k, j = k \end{cases}$$

$$\mathbf{fl} (\eta_k^k \ \mathbf{m}_{k,j}^k) & i = k, j > k$$

$$\mathbf{fl} (\mathbf{m}_{i,j}^k + \eta_i^k \ \mathbf{m}_{k,j}^k) & i > k, j > k$$

$$\mathbf{m}_{i,j}^k & \text{otherwise}$$
 (2)

i.e.,

Thus if

$$\mathbf{L}^{k} = \begin{bmatrix} \mathbf{1} & & & & \\ & \mathbf{1} & & & & \\ & & \ddots & & & \\ & & & \mathbf{1} & \mathbf{1} & \\ & & & \mathbf{1} & \mathbf{1} & \\ & & & & & & & & \mathbf{1} & \\ & & & & &$$

and E^{k} is the matrix with elements ϵ_{ij}^{k} ,

then

$$M^{k+1} = L^k M^k + E^k$$
 (4)

$$\epsilon_{ij}^{k} = \begin{cases} m_{kk}^{k} \eta_{k}^{k} \delta_{kk}^{k} & i = j = k \\ m_{ik}^{k} \delta_{ik}^{k} + m_{kk}^{k} \eta_{i}^{k} \delta_{kk}^{k} & i > k, j = k \\ m_{kj}^{k} \eta_{k}^{k} \delta_{kj}^{k} & i = k, j > k \\ m_{ij}^{k} \delta_{ij}^{k} + 2 m_{kj}^{k} \eta_{i}^{k} \delta_{ij}^{k} & i > k, j > k \\ 0 & \text{otherwise} \end{cases}$$

$$(5)$$

and
$$|\delta_{ij}^k| \le \epsilon_1^i$$
, $|\delta_{ij}^{ik}| \le \epsilon_1^i$, for all i, j, k.

After s iterations $M^{S+1} = U$, (say), an upper triangular matrix which has ones down the diagonal. Let

$$F^{k} = (L^{k})^{-1} E^{k}$$

$$M^{k+1} = L^{k} M^{k} + L^{k} F^{k}$$

$$= L^{k} (M^{k} + F^{k})$$

But the elements of F^k are zero except for $i \ge k$ and $j \ge k$ (because those of E^k are).

...
$$L^{i} F^{k} = F^{k}$$
 for $i < k$
... $L^{k-1} L^{k-2} ... L^{2} L^{1} F^{k} = F^{k}$

$$U = M^{S+1}$$

$$= L^{S} L^{S-1} \cdot \cdot \cdot L^{2} L^{1} M^{1} + \sum_{k=1}^{S} (L^{S} L^{S-1} \cdot \cdot \cdot L^{k} F^{k})$$

$$= L^{S} L^{S-1} \cdot \cdot \cdot L^{2} L^{1} M^{1} + \sum_{k=1}^{S} (L^{S} L^{S-1} \cdot \cdot \cdot L^{2} L^{1} F^{k})$$

$$= L^{S} L^{S-1} \cdot \cdot \cdot L^{2} L^{1} (M + F) \cdot \cdot \cdot \cdot (6)$$

$$F = \sum_{k=1}^{s} F^{k}.$$

If

$$U = \begin{bmatrix} 1 & u_{12} & u_{13} & \cdots & u_{1s} \\ & 1 & u_{23} & & u_{2s} \\ & & & \ddots & & u_{3s} \\ & & & \ddots & & \vdots \\ & & & & 1 \end{bmatrix}$$

then

$$u^{-1} = v^2 u^3 \dots u^{s-1} u^s$$

from which

$$U^2 U^3 \dots U^{s-1} U^s L^s L^{s-1} \dots L^2 L^1 (M+F) = I.$$

A bound on ||F|| is now required. Since

$$F^{k} = (L^{k})^{-1} E^{k},$$

$$F_{ij}^{k} = \begin{cases} m_{ik}^{k} \delta_{ik}^{k} & i \geq k, j = k \\ m_{ij}^{k} \delta_{kj}^{k} & i = k, j > k \\ m_{ij}^{k} \delta_{ij}^{k} + 2 \eta_{i}^{k} m_{kj}^{k} (\delta_{ij}^{k} - \delta_{kj}^{k}) & i > k, j > k \\ 0 & \text{otherwise} \end{cases}$$

$$(7)$$

An important point is that these error terms do not involve the 'pivot' element η_k^k Summing over k gives

$$F_{ij} = \begin{cases} m_{ij}^{j} \delta_{ij}^{j} + \sum_{k=1}^{j-1} (m_{ij}^{k} \delta_{ij}^{k} + 2 \eta_{i}^{k} m_{kj}^{k} (\delta_{ij}^{k} - \delta_{kj}^{k})) \\ i \geq j \end{cases}$$

$$m_{ij}^{i} \delta_{ij}^{i} + \sum_{k=1}^{i-1} (m_{ij}^{k} \delta_{ij}^{k} + 2 \eta_{i}^{k} m_{kj}^{k} (\delta_{ij}^{k} - \delta_{kj}^{k}))$$

$$i < j \qquad (8)$$

Let
$$\rho = \max_{i,j,k} \{ | m_{ij}^k | : i \ge k \}$$

$$M = \max_{i,k} \{ | m_{ij}^k | : i \ge k \}$$

$$q_i = \sum_{k=1}^{s} \nabla (n_i^k)$$

$$\sigma_{ij} = \sum_{k=1}^{s} \nabla (m_{ij}^k)$$

Then
$$|F_{ij}| \le (\sigma_{ij} + 4q_i) \rho M \epsilon'_i$$

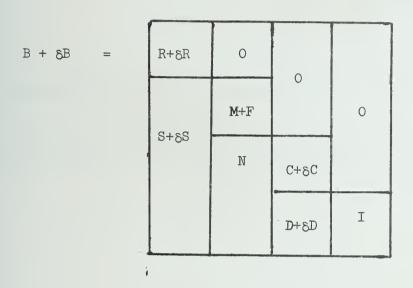
 $\therefore ||F|| \le s (\sigma + 4q) \rho M \epsilon'_i$

where
$$q = \max_{i} (q_i)$$

$$\sigma = \max_{i,j} (\sigma_{ij}).$$

This gives an a posteriori bound for ||F|| in terms of the sparseness and the size of the elements of the inverse and the intermediate matrices M^1 , M^2 , . . . , M^p . As might be expected, sparse matrices give much smaller errors.

Combining these results indicates that the computed inverse is the exact inverse, in product form, of the perturbed matrix



It is clear from equation (9) that both sparseness and size of the elements are crucial factors in bounding the error. The previous chapter dealt with preserving sparseness. Consideration now must be given to bounding the size of the elements in the Π vectors and the intermediate matrices that are formed. In forming an Π vector, start with some column vector $v = [v_1, v_2, \dots, v_p]^T$ (say) and 'pivot' on some element v_r , giving

$$\eta = \text{fl} \left[-v_1/v_r, \dots, -v_{r-1}/v_r, 1/v_r, -v_{r+1}/v_r, \dots, -v_p/v_r \right]^T$$

where the value of r may be arbitrarily chosen. The usual method is to select the pivot element so that $|\mathbf{v}_r|$ is greater than some 'tolerance' value so that the elements of η will not be 'too large'. Since, as equation (7) shows, the element $1/\mathbf{v}_r$ does not appear in any error term, the remaining terms may

be bounded by c (say) by insisting that r be chosen such that:

$$|v_r| \ge \frac{1}{c} |v_i|$$
 for $i = 1, ..., r-1$, $r+1, ..., p$

It may be noted that for dense systems, one normally selects the largest element of v for this pivot, which is equivalent to taking c = 1. Thus the above scheme is a compromise between minimizing the number of nonzero elements and reducing their size.

7. ITERATIVE REFINEMENT

Suppose that an approximate inverse to the basis B and an approximate solution vector x have been calculated. It is desired to calculate a more accurate solution using iterative refinement. A detailed analysis of iterative refinement using floating point arithmetic is given by Moler [7]. A simple adaption of that work is developed here to cover the sparse matrix case.

Each iteration of the refinement requires the following three steps:

- 1) Compute the residuals: $r^m = b Bx^m$ with $x^1 = x$.
- 2) Compute the correction: $d^{m} = B^{-1}r^{m}$.
- 3) Add the correction to form: $x^{m+1} = x^m + d^m$.

The step where most accuracy is required is in calculating the residuals. It is assumed that extended precision is used in this step. Computing $d^{m} = B^{-1}r^{m}$ is done using the approximate inverse to B.

Let

$$r^{m} = b - B x^{m} = c^{m}$$
 (1)

$$d^{m} = (B + E^{m})^{-1} r^{m}$$
 (2)

$$x^{m+1} = x^m + a^m + g^m \tag{3}$$

where c^m , E^m and g^m are correction terms due to round-off error and are introduced to make the equations hold exactly. A bound for $||E^m||$ is known, using the results of either Chapter 4 or Chapter 6. It is required that the inverse be reasonably good so that

$$||E^{m}|| < (||B^{-1}||)^{-1}$$

Let

$$F = B^{-1} E^{m}$$

Then

$$||F^{m}|| = ||B^{-1}E^{m}|| \le ||B^{-1}|| ||E^{m}|| < 1.$$

and hence $I + F^{m}$ is nonsingular.

Equation (2) becomes

$$d^{m} = (I + F^{m})^{-1} B^{-1} r^{m}.$$
 (4)

Suppose x is the exact solution $x = B^{-1} b$.

Then

$$x^{m+1} - x = x^m - x + d^m + g^m$$

$$= (I+F^m)^{-1} [(I+F^m) (x^m-x) + B^{-1}(b-Bx^m + c^m)] + g^m$$

$$= (I+F^m)^{-1} [F^m(x^m-x) + B^{-1} c^m] + g^m$$

Hence

$$||x^{m+1} - x|| \le \frac{||F^m||}{1 - ||F^m||} ||x^m - x|| + \frac{||B^{-1}||}{1 - ||F^m||} ||c^m|| + ||g^m||$$

Suppose an upper bound for $|E^{m}|$ is given by $|\delta B|$. Then

$$||x^{m+1} - x|| \le \frac{||B^{-1}||}{1 - ||\delta B|| ||B^{-1}||} (||\delta B|| ||x^m - x|| + ||c^m||) + ||g^m||$$

Estimates of $||c^{m}||$ and $||g^{m}||$ are now required.

$$g^{m} = fl(x^{m} + d^{m}) - (x^{m} + d^{m})$$

$$\vdots \quad \mathbf{g}_{\mathbf{i}}^{\mathbf{m}} = \boldsymbol{\theta}_{\mathbf{i}} \mathbf{x}_{\mathbf{i}}^{\mathbf{m}} + \boldsymbol{\theta}_{\mathbf{i}}^{\mathbf{i}} \mathbf{d}_{\mathbf{i}}^{\mathbf{m}} \quad , \quad \text{for } |\boldsymbol{\theta}_{\mathbf{i}}| \leq \boldsymbol{\epsilon}_{\mathbf{l}}, |\boldsymbol{\theta}_{\mathbf{i}}^{\mathbf{i}}| \leq \boldsymbol{\epsilon}_{\mathbf{l}}$$

$$||g^{m}|| \le \epsilon_{1} (||x^{m}|| + ||d^{m}||)$$

When the process converges, $||d^{m}|| \le ||x^{m}||$

$$\dots \quad ||g^m|| \leq 2 \, \varepsilon_1 \, ||x^m||$$

$$\leq 2 \in_{1} (||x|| + ||x^{m} - x||).$$

In calculating r^m , it is assumed that each component $r^m_i = \operatorname{fl}_2(b_i - \sum b_{ij} x^m_j)$ is calculated using extended precision and the result then is rounded to normal precision. This is denoted by $\operatorname{fl}_{2,1}(b_i - \sum b_{ij} x^m_j)$.

$$c^{m} = r^{m} - b + Bx^{m}$$

and

$$c_{i}^{m} = fl_{2,1} (b_{i} - \sum B_{ij} x_{j}^{m}) - (b_{i} - \sum B_{ij} x_{j}^{m})$$

$$= (fl_{2} (b_{i} - \sum B_{ij} x_{j}^{m})) (1 + \theta_{i}) - (b_{i} - \sum B_{ij} x_{j}^{m})$$

$$= \{b_{i}(1 + \xi_{i}) - \sum B_{ij} x_{j}^{m} (1 + \gamma_{ij})\} (1 + \theta_{i}) - (b_{i} - \sum B_{ij} x_{j}^{m})$$

$$= (b_{i} - \sum B_{ij} x_{j}^{m}) \theta_{i}$$

$$+ \{b_{i} \xi_{i} - \sum B_{ij} x_{j}^{m} \gamma_{ij}\} (1 + \theta_{i})$$

where
$$|\theta_{i}| \leq \epsilon_{1}, |\xi_{i}| \leq \epsilon_{2}$$

$$|\gamma_{ij}| \leq (q+2) \epsilon_{2}^{i}$$

$$q = \max_{i} \sum_{j=1}^{m} \nabla (B_{ij})$$

where $k(B) = ||B|| ||B^{-1}||$, the condition number of B.

$$= \sigma_1 ||x^m - x|| + \sigma_2 ||x||$$
 (say).

$$\cdot \cdot \frac{||\mathbf{x}^{m+1} - \mathbf{x}||}{||\mathbf{x}||} \leq \sigma_1 \frac{||\mathbf{x}^m - \mathbf{x}||}{||\mathbf{x}||} + \sigma_2$$

Since x^{1} is calculated in the same way (using, in effect, $x^{0} = 0$)

$$||x^{1} - x|| \le \sigma_{1} ||x||$$

(This may be shown by setting $\epsilon_2' = 0$ for this step (no calculation of residuals is involved) and noting that the only other term in σ_2 , $2\epsilon_1$, came from $2\epsilon_1 ||x^m|| \le 2\epsilon (||x|| + ||x^m - x||)$, but $x^0 = 0$.

$$\cdots \frac{||x||}{||x_1 - x||} \leq a^{\frac{1}{2}}$$

... by induction,

$$\frac{\left|\left|\left|x^{m}-x\right|\right|}{\left|\left|x\right|\right|} \leq \frac{\sigma_{1}^{m}+\sigma_{2}}{1-\sigma_{1}}, \quad \sigma_{1} \leq 1.$$

Therefore, provided $\sigma_1 < 1$, convergence of $\mathbf{x}^{\mathbf{m}}$ to within the tolerance given by

$$\frac{\sigma_2}{1-\sigma_1} \quad ||x||,$$

can be guaranteed.

depends essentially on $||\delta B||$, $||B^{-1}||$ and the measure of sparseness, q. It is essential that B be well conditioned and that $||\delta B||$ be small, i.e., an accurate inverse of B is known.

The effect of a small value of q is to reduce the need for using double precision in calculating the residuals, since it is always multiplied by ϵ_2^{\prime} .

8. CONCLUDING REMARKS

The question remains as to what constitutes an optimal feasible solution. Strictly speaking, it is one for which $x \ge 0$ and the reduced costs d_j are also ≥ 0 for j=1(1)n. In practice, small negative tolerances, δ_f , δ_o (say) are chosen so that the solution is considered optimal and feasible provided

$$x_{i} \geq \delta_{f}$$
 for all i

$$d_j \ge \delta_o$$
 for all columns j not in the basis.

An alternative approach is as follows: Determine some sufficiently refined solution x and a basis B, such that

$$Bx = b$$

may be considered exact. If some x_i are < 0, define

$$\bar{x} = x + \delta x \ge 0$$

where

$$\delta x_{i} = \begin{cases} 0 & , & x_{i} \geq 0 \\ -x_{i} & , & x_{i} < 0 \end{cases}$$

This produces a change in b, 8b (say), where

$$\delta b = B \delta x$$

The vector &b can either be calculated directly, or its norm may be estimated from

If the elements of δb are sufficiently small that $b+\delta b$ is an acceptable right hand side, then \bar{x} may be considered to be the required solution to the problem.

Similarly, in calculating $d_j = c_j - \sum_{i=1}^m \pi_i \ a_{ij}$, π may be refined in the same way as x, so that the d_j are accurately known. If $d_j < 0$ for some j, a small change in the corresponding c_j will set them to zero. If this new cost vector is acceptable, then the problem is solved.

The techniques described in Chapters 3 through 6 and above enable the exact solution to the slightly perturbed problem:

minimize

$$z + \delta z = (c + \delta c)^{T} x$$

subject to

$$(A + \delta A) x = b + \delta b$$

 $x > 0.$

to be obtained. If $||\delta A||$, $||\delta b||$ and $||\delta c||$ are within acceptable limits, then the problem may be considered solved. If more accuracy is required, iterative refinement as described in Chapter 7 may be used.

LIST OF REFERENCES

- [1] Bartels, Richard H., "A Numerical Investigation of the Simplex Method", Technical report no. CS 104, Stanford University, 1968.
- [2] Clasen, R. J., "Techniques for Automatic Tolerance Control in Linear Programming", Comm. ACM, 2 (November 1966), pp. 802-803.
- [3] Dantzig, George B., Linear Programming and Extensions, Princeton University Press, 1965.
- [4] Dantzig, George B., Harvey, Roy P., McKnight, Robert D., Smith, Stanley S.,
 "Sparse Matrix Techniques in two Mathematical Programming Codes",
 Technical report no. 69-1, Stanford University, 1969.
- [5] Forsythe, George E. and Moler, Cleve B., Computer Solution of Linear Algebraic Systems, Prentice Hall, 1967.
- [6] Hadley, G., Linear Programming, Addison Wesley, 1962.
- [7] Moler, Cleve B., "Iterative Refinement in Floating Point", JACM, 14 (April 1967), pp. 316-321.
- [8] Orchard-Hays, William, Advanced Linear Programming Techniques, McGraw Hill, 1968.
- [9] Wilkinson, J. H., Rounding Errors in Algebraic Processes, Prentice Hall, 1963.

APPENDIX

Lemma 1: The solution of a triangular system of equations.

Suppose that the components of a vector x are computed in the order x_1, x_2, \dots, x_n by

$$x_1 = b_1$$

 $x_i = fl (b_i + \sum_{k=1}^{i-1} r_{ik} x_k)$ $i = 2, 3, ..., r$

or, equivalently, by the set of equations

i.e., Rx = b, is solved. Then x is the exact solution of the perturbed triangular system:

$$(R + \delta R) x = b + \delta b,$$

where,

$$\delta R = \begin{bmatrix} 0 & & & & & & \\ -r_{21}\delta_{21} & & 0 & & & \\ -r_{31}\delta_{31} & -r_{32}\delta_{32} & & 0 & & \\ \vdots & & & & \ddots & \\ -r_{ml}\delta_{ml} & & & -r_{mm-l}\delta_{mm-l} & 0 \end{bmatrix}$$

and

$$\delta b = [b_1 \delta'_1, \dots, b_m \delta'_m]^T,$$

$$|\delta'_i| \leq q_i \epsilon'_i,$$

$$|\delta_{ik}| \leq (q_i + 1) \epsilon'_i$$

with

$$q_{i} = \sum_{k=1}^{i-1} \nabla (r_{ik})$$

Proof

In Chapter 2, it was shown that

fl
$$(\sum_{i=1}^{n} x_{i}y_{i}) = \sum_{i=1}^{n} x_{i}y_{i}$$
 (1+8_i)

where

$$\mid \delta_i \mid \leq (q_x + 1) \in '$$

and

$$q_{x} = \sum_{i=1}^{n} \nabla (x_{i})$$

Similarly, if $p = fl(z + \sum_{i=1}^{n} x_i y_i)$ is calculated from:

$$s_1 = z$$

$$p = s_{n+1}$$

then

fl
$$(z + \sum_{i=1}^{n} x_{i}y_{i}) = z (1 + \Delta) + \sum_{i=1}^{n} x_{i}y_{i} (1 + \delta_{i})$$

where

$$|\Delta| \le q_x \in '$$
,

$$\mid \delta_{\mathtt{i}} \mid \leq (q_{\underline{x}} + 1) \in \mathsf{'}$$

and

$$q_{x} = \sum_{i=1}^{n} \nabla (x_{i})$$
.

Using this result,

$$x_i = fl (b_i + \sum_{k=1}^{i-1} r_{ik} x_k)$$

=
$$b_{i} (1 + \delta_{i}) + \sum_{k=1}^{i-1} r_{ik} x_{k} (1 + \delta_{ik})$$

where

$$|\delta_{ik}| \leq (q_i + 1) \epsilon_i'$$

and

$$\mid \delta_{i}^{'} \mid \leq q_{i} \in 1$$
.

Hence

$$x_i = b_i + b_i \delta_i' + \sum_{k=1}^{i-1} (r_{ik} + r_{ik} \delta_{ik}) x_k$$

which completes the proof.

Note also that

$$||\delta R|| = \max_{i} \sum_{j=1}^{i-1} |r_{ij} \delta_{ij}|$$

$$\leq \max_{i} (q_{i} \max_{1 \leq j \leq i-1} |r_{ij} \delta_{ij}|)$$

$$\leq q (q+1) \epsilon'_{1 i,j} |r_{ij}|$$

where $q = \max_{i} q_{i}$.

If r_{ik} results from a division, as described earlier, two extra errors are involved so that, in the lemma,

$$|\delta_{ik}| \leq (q_i + 3) \epsilon_i'$$

and $|8R| \le q (q + 3) \in \max_{i,j} |r_{ij}|$

Lemma 2. A direct error bound on x.

Under the same hypotheses as in Lemma 1, if \bar{x} is the exact solution $R^{-1}b$, and $\delta x=x-\bar{x}$, then

$$||\delta x|| \le \frac{||R^{-1}||}{1-2m||R^{-1}||} [m||x|| + \epsilon_1' ||b||]$$

where

$$m = q (q + 2) \in \max_{i,k} |r_{ik}|,$$

provided that

$$m ||R^{-1}|| < \frac{1}{2}$$
.

Proof

As before,

$$(R+\delta R) x = b + \delta b$$

also

$$R \bar{x} = b$$

 $b + 8b = (R + 8R) (\bar{x} - 8x)$

$$= R\bar{x} + 8R\bar{x} - (R + 8R) 8x$$

$$...$$
 $\delta b = \delta R \bar{x} - (R + \delta R) \delta x$

$$8x = (R + 8R)^{-1} [8R\bar{x} - 8b]$$

and
$$||\delta x|| \le ||(R + \delta R)^{-1}||[||\delta R||||\bar{x}|| + ||\delta b||]$$

but
$$||(R + \delta R)^{-1}|| = ||[R(I + R^{-1}\delta R)]^{-1}||$$

$$\leq ||R^{-1}|| ||(I + R^{-1} \delta R)^{-1}||$$

$$\leq \frac{||R^{-1}||}{1-||R^{-1}\delta R||}$$

$$\leq \frac{||R^{-1}||}{1-||R^{-1}|| ||8R||}$$

provided $||R^{-1}|| ||8R|| < 1$

(and hence $||R^{-1}\delta R|| < 1$)

But $\bar{x} = x - \delta x$,

$$||\bar{x}|| \le ||x|| + ||\delta x||$$

which reduces to

$$||\delta x|| \le \frac{||R^{-1}||}{1-2||R^{-1}|| ||\delta R||} [||\delta R|| ||x|| + \epsilon_1' ||b||]$$

provided

$$||R^{-1}|| ||\delta R|| < \frac{1}{2}$$

and
$$m = |\delta R|$$
.



Security Classification

DOCUMENT	CONTRACT	PATA	D 8	179
PULUMENI	CUNTRUL	UAIA -	K 6	w

DOCUMENT CONT	ROL DATA - R	& D			
(Security classification of title, body of abstract and indexing	emotetten must be	entered when the	overall report is ctassified)		
Department of Computer Science		UNCLASSIFIED			
Urbana, Illinois 61801					
REPORT TITLE					
AN ERROR ANALYSIS OF SOLUTIONS TO SPARSE	LINEAR PRO	GRAMMING PI	ROBLEMS		
DESCRIPTIVE NOTES (Type of report and inclusive dates)					
Research Report					
AUTHOR(S) (First name, middle initial, last name)					
Raymond J. Lermit					
REPORT DATE	76. TOTAL NO.	NE PAGES	76. NO. OF REFS		
		,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,			
September 17, 1969	69	'S REPORT NUME	9		
	- Oktobal Con	The rest is the second	2010)		
46-26-15-305	Dag Dane				
	DUS Repo	rt No. 351			
US AF 30(602)4144	90. OTHER REPO	RT NO(5) (Any of	her numbers that may be as	alaned	
	this report)	, ,			
DISTRIBUTION STATEMENT					
Qualified requesters may obtain copies fr	com DCS.				
	. 0.1.1				
SUPPLEMENTARY NOTES		MILITARY ACTI			
	4	r Developme			
NONE		s Air Force			
	Rome, No	ew York 13 ¹	+40		

This thesis discusses the round-off errors incurred in solving Linear Programming problems by the Revised Simplex Method, using the Product Form of the Inverse - the most usual method of solution. The matrices involved are assumed to be sparse as is usually the case.

Also discussed are the reinversion of the basis matrix and iterative refinement of the solution obtained.

. ABSTRACT

UNCLASSIFIED

Security Classification		LINKA		LHK		LINKC	
14 KEY WORDS	ROLE	WT	ROLE	WT	ROLE	WT	
	,						
Linear programming							
Error analysis							
Sparse matrices							
Revised product method							
Matrix inversion							
						-	

UNCLASSIFIED
Security Classification



















UNIVERSITY OF ILLINOIS-URBANA 510 84 IL6R no. C002 no.349-354(1969 Report /

